

# Stein's Method

Matthew Kwan

November 1, 2013

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Notation . . . . .	2
<b>I</b>	<b>Theory</b>	<b>2</b>
<b>2</b>	<b>General Probability Theory</b>	<b>3</b>
2.1	Review of basic concepts . . . . .	3
2.2	Coupling . . . . .	7
2.3	Markov Chains . . . . .	7
2.4	The Weak Topology on Probability Measures . . . . .	8
<b>3</b>	<b>Random Combinatorial Structures</b>	<b>13</b>
<b>4</b>	<b>Stein's Method in the Abstract</b>	<b>14</b>
4.1	The method of Exchangeable Pairs . . . . .	19
4.2	Size-Bias Coupling . . . . .	26
<b>II</b>	<b>Applications</b>	<b>26</b>

## 1 Introduction

**Comment 1.1.** The plan is to introduce with limit theorems: Central Limit theorem, Poisson Limit theorem. The failure of limit theorems is that they provide no understanding of speed of convergence, in particular convergence cannot be assumed to be uniform as parameters vary.

Stein’s method is a technique for bounding the distance between distributions, with a variety of different distance metrics. Quantitative bounds can be useful in their own right, or can be further applied to prove asymptotic results.

## 1.1 Notation

For this thesis, the set of natural numbers  $\mathbb{N}$  includes zero. We write  $1_A$  for the characteristic function of  $A$ :  $1_A(x) = 1$  if  $x \in A$ , otherwise  $1_A(x) = 0$ . Also,  $[k]$  denotes the set  $\{1, \dots, k\}$ .

Unless otherwise specified, all asymptotics are as  $n \rightarrow \infty$ . Apart from standard asymptotic notation, we use two notions of asymptotic equivalence:  $f \sim g$  means  $f = g(1 + o(1))$  and  $f \asymp g$  means  $f = O(g)$  and  $g = O(f)$ .

In this thesis, unless stated otherwise, graphs are labelled. That is, they are distinguished even within isomorphism classes. A graph may not have loops or multiple edges; an object which is allowed to have loops and/or multiple edges will be called a multigraph.

The phrase “randomly choose” is taken to mean a uniformly random choice. That means, each possible option is chosen with equal probability.

## Part I

# Theory

## 2 General Probability Theory

### 2.1 Review of basic concepts

**Comment 2.1.** I'm a little bit uncertain how much depth to go into for this. At the moment, it's written so that someone who's seen measure theory but no probability theory (an analyst) can understand. Where possible, I've tried to translate things into the discrete case, because it's often more intuitive (and since I plan for applications to be combinatorial).

For many combinatorial applications, an informal understanding of probability theory will suffice. However, in this thesis a rigorous foundation in probability theory will be useful. The following is intended only as a brief review.

**Definition 2.1.** A *probability space* is a measure space  $(\Omega, \mathcal{A}, \mathbb{P})$  with  $\mathbb{P}(\Omega) = 1$ . In this case we say  $\mathbb{P}$  is a *probability measure*, and denote the set of all probability measures on  $(\Omega, \mathcal{A})$  by  $\mathcal{P}(\Omega, \mathcal{A})$  or  $\mathcal{P}(\Omega)$  if there is no ambiguity. An *event* is a measurable set  $A \in \mathcal{A}$ . When we have only defined one measure on the set  $\Omega$ , we will sometimes abuse notation and write  $\mathbb{P}$  to mean  $\text{mean}(\Omega, \mathcal{A}, \mathbb{P})$ , and vice versa.

**Definition 2.2.** A *probability space* is a measure space  $(\Omega, \mathcal{A}, \mathbb{P})$  with  $\mathbb{P}(\Omega) = 1$ . In this case we say  $\mathbb{P}$  is a *probability measure*, and denote the set of all probability measures on  $(\Omega, \mathcal{A})$  by  $\mathcal{P}(\Omega, \mathcal{A})$  or  $\mathcal{P}(\Omega)$  if there is no ambiguity. An *event* is a measurable set  $A \in \mathcal{A}$ .

For our purposes  $\Omega$  will often be a finite set of combinatorial objects, with  $\mathcal{A}$  as the power set of  $\Omega$ . In this case  $\mathbb{P}$  is defined by  $\mathbb{P}(\omega) := \mathbb{P}(\{\omega\})$ , for each  $\omega \in \Omega$ . We will discuss specific probability spaces on combinatorial objects in Section 3.

For an event  $A$ ,  $\mathbb{P}(A)$  is interpreted as the “probability that  $A$  occurs”. For combinatorial spaces, events are usually of the form  $A = \{\omega \in \Omega : P(\omega) \text{ holds}\}$ , where  $P(\omega)$  is some property of an object  $\omega$ . For clarity, we often abuse notation slightly and write  $\mathbb{P}(P(\omega) \text{ holds})$  instead of  $\mathbb{P}(A)$ .

**Definition 2.3.** A *random element*  $X : \Omega_1 \rightarrow \Omega_2$  is a measurable function from a probability space  $(\Omega_1, \mathcal{A}_1, \mathbb{P})$  to some measure space  $(\Omega_2, \mathcal{A}_2, \mu)$ . If the target measure space is  $\mathbb{R}^n$  with the Borel  $\sigma$ -algebra and the Lebesgue measure, then we say  $X$  is a *random vector*; if  $n = 1$  then  $X$  is a *random variable*. If  $X$  only takes countably many values then we say  $X$  is *discrete*.

If the underlying probability space  $\Omega_1$  is countable, then any function is measurable.

We will often be interested in the probability that a random element takes certain values, without regard to the underlying probability space.

**Definition 2.4.** Suppose  $X$  is a random element with target measure space  $(\Omega, \mathcal{A}, \mu)$ . The *distribution* (or *law*)  $\mathcal{L}_X$  of  $X$  is the pushforward measure with respect to  $X$ . That is, it is a probability measure defined by  $\mathcal{L}_X(A) = \mathbb{P}(X^{-1}(A))$  for  $A \subseteq \mathcal{A}$ . We also occasionally use the notation  $\mathcal{L}(X) := \mathcal{L}_X$  for ease of reading.

It is worth noting that in fact any probability measure is the distribution of some random element. To see this, note that given a probability measure  $\mathbb{P} \in \mathcal{P}(\Omega)$ , we can choose  $X = \text{id}_\Omega$  to have  $\mathcal{L}_X = \mathbb{P}$ . So, it is often convenient to specify random variables by their distributions, without defining an underlying probability space.

**Definition 2.5.** The notation  $X \in \mathcal{L}$  means  $\mathcal{L}_X = \mathcal{L}$ .

We can use slightly abusive (but standard) notation like  $\mathbb{P}(X > 1)$  to denote  $\mathcal{L}_X(\{x : x > 1\})$ . This is equal to  $\mathbb{P}(\{\omega \in \Omega : X(\omega) > 1\})$  for any particular realization of  $X$  as a function on a probability space  $(\Omega, \mathcal{A}, \mathbb{P})$ .

**Definition 2.6.** The *support*  $\text{supp}(X)$  of a discrete random element  $X$  is the set  $\{k \in \Omega : \mathbb{P}(X = k) > 0\}$ .

**Example 2.7.** If  $X$  has the normal distribution with parameters  $\mu$  and  $\sigma$  then we say  $\mathcal{L}_X = \mathcal{N}(\mu, \sigma)$ ; this distribution is defined by  $\mathcal{L}_X(B) = \frac{1}{\sigma\sqrt{2\pi}} \int_B e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$  for Borel  $B$ .

If  $X$  is discrete, then  $\mathcal{L}_X$  is just an assignment of a probability to each possible value.

**Example 2.8.** If  $X$  is Poisson distributed with parameter  $\lambda$ , we write  $\mathcal{L}_X = \text{Po}(\lambda) = \text{Po}_\lambda$ ; this is defined by  $\mathbb{P}(X = k) = \frac{\lambda^k e^{-\lambda}}{k!}$ .

**Definition 2.9.** The *expected value* of a random variable  $X$  is  $\mathbb{E}X = \int x \, d\mathcal{L}_X(x)$ .

For a random variable  $X$  that takes integer values, this definition is equivalent to the well-known formula  $\mathbb{E}X = \sum_{x \in \mathbb{Z}} x \mathbb{P}(X = x)$ .

*Remark 2.10.* If  $X$  is a random variable that can be interpreted as counting the number of objects that satisfy some property, then we can express  $X$  as a sum of indicator variables  $\sum_i 1_{A_i}$ , where  $A_i$  is the event that the  $i$ th object satisfies our property. We have  $\mathbb{E}X = \sum_i \mathbb{E}1_{A_i} = \sum_i \mathbb{P}A_i$ . So, in order to compute the expectation of  $X$  we just need to compute the probability that each object satisfies our required property.

If we fix a particular underlying probability space  $(\Omega, \mathcal{A}, \mathbb{P})$ , we can also equivalently view expectation as a linear functional on the space of integrable functions:  $\mathbb{E}X = \int X(\omega) \, d\mathbb{P}$ . So,  $\mathbb{E}$  is defined in terms of a particular underlying probability space. Sometimes we will define a new probability space  $(\Omega, \mathcal{A}, \mathbb{P}')$  by changing the measure on the same underlying set. In this case we will write  $\mathbb{E}_{\mathbb{P}'}$  to indicate expectation with respect to the measure  $\mathbb{P}'$ , to avoid ambiguity.

In fact, the expectation functional defines its underlying probability measure, because  $\mathbb{E}1_A = \mathbb{P}(A)$ . Since the distribution of a random variable is specified by a probability measure, the distribution  $\mathcal{L}(X)$  of a random variable  $X$  also uniquely defines an expectation functional  $\mathbb{E}_X := \mathbb{E}_{\mathcal{L}(X)}$ .

**Definition 2.11.** For two collections  $S, S' \subseteq \mathcal{A}_1$  of events, we say that  $S$  and  $S'$  are *independent* if  $\mathbb{P}(A \cap A') = \mathbb{P}(A)\mathbb{P}(A')$  for each  $A \in S$  and  $A' \in S'$ . If  $S = \{A\}$  contained a single set, then we say  $A$  itself is independent of  $S'$ .

**Definition 2.12.** Let  $(\Omega_1, \mathcal{A}_1, \mathbb{P})$  be a probability space and  $(\Omega_2, \mathcal{A}_2, \mu)$  a measure space. Let  $X$  be a random variable  $\Omega_1 \rightarrow \Omega_2$ , and let  $S$  be the set of all events of the form  $\{\omega \in \Omega_1 : X(\omega) \in A_2\}$  for  $A_2 \in \mathcal{A}_2$ . If  $S$  is independent of  $S'$  then we say  $X$  itself is independent of  $S'$ .

We can analogously say that two random variables are independent, or a random variable and an event are independent, or any similar combination.

**Definition 2.13.** If two objects are not independent, then we say they are *dependent*.

**Definition 2.14.** Suppose  $X : \Omega_1 \rightarrow \Omega_2$  is a random element defined on these spaces, and  $A_1 \in \mathcal{A}_1$  is an event with nonzero probability. Then the *distribution of  $X$  conditioned on  $A_1$*  is denoted by  $\mathcal{L}_{X|A_1}$  and defined by  $\mathcal{L}_{X|A_1}(A_2) = \mathbb{P}(X \in A_2|A_1)$  for  $A_2 \in \mathcal{A}_2$ . The expected value of a random variable with distribution  $\mathcal{L}_{X|A_1}$  is called the *conditional expected value of  $X$  given  $A_1$*  and is denoted  $\mathbb{E}[X|A_1]$ .

We can also define conditional expectation with respect to another random variable. If  $X_1$  and  $X_2$  are random variables defined on the same underlying probability space  $(\Omega, \mathcal{A}, \mathbb{P})$ , then the sets  $X_2^{-1}(B)$  for Borel  $B$  comprise a sub- $\sigma$ -algebra  $\mathcal{A}'$  of  $\mathcal{A}$ . Then,  $\mu : A' \mapsto \mathbb{E}[X_1 1_{A'}]$  is a signed measure on  $\mathcal{A}'$  that is absolutely continuous with respect to the restriction of  $\mathbb{P}$  to  $\mathcal{A}'$ . By the Radon-Nikodym theorem there is an  $\mathcal{A}'$ -measurable random variable  $\mathbb{E}[X_1|X_2]$  that satisfies  $\mathbb{E}[X_1 1_{A'}] = \mathbb{E}[\mathbb{E}[X_1|X_2] 1_{A'}]$  for all  $A'$  in  $\mathcal{A}'$ . This random variable is almost uniquely defined: for any two choices of  $\mathbb{E}[X_1|X_2]$ , the probability that they differ is zero.

**Definition 2.15.** The random variable  $\mathbb{E}[X_1|X_2]$  as defined above is called the *conditional expectation of  $X_1$  with respect to  $X_2$* . We can also view conditional expectation as a linear operator between functions: we define  $\mathbb{E}^{X_2}$  by  $X_1 \mapsto \mathbb{E}[X_1|X_2]$ .

This definition generalizes the previous definition of expectation conditioned on an event: if  $\omega \in A$  and  $\mathbb{P}(A) > 0$  then  $\mathbb{E}[X 1_A](\omega) = \mathbb{E}[X|A]$ .

Note that if  $X_2$  is discrete then we do not need to invoke Radon-Nikodym. We can define  $\mathbb{E}[X_1|X_2]$  by  $\mathbb{E}[X_1|X_2](\omega) = \mathbb{E}[X_1|X_2 = X_2(\omega)]$  whenever  $\mathbb{P}(X_2 = X_2(\omega)) > 0$ .

We finally present a simple consequence of the definition of conditional expectation.

**Proposition 2.16** (Tower Law of Expectation).  $\mathbb{E}[\mathbb{E}^{X_2} X_1] = \mathbb{E}[X_1]$

*Proof.*  $\mathbb{E}[\mathbb{E}^{X_2} X_1] = \mathbb{E}[\mathbb{E}[X_1|X_2] 1_\Omega] = \mathbb{E}[X_1 1_\Omega] = \mathbb{E}[X_1]$  □

## 2.2 Coupling

Given a finite collection of measure spaces  $(\Omega_1, \mathcal{A}_1, \mu_1), \dots, (\Omega_n, \mathcal{A}_n, \mu_n)$  recall the construction of the product measure space  $(\Omega, \mathcal{A}, \mu) := (\prod_{i=1}^n \Omega_i, \otimes_{i=1}^n \mathcal{A}_i, \prod_{i=1}^n \mu_i)$ . If a random element takes values in a product space then each component is measurable, and conversely if the components of a random tuple are measurable then that tuple is measurable in the product space. So, we can make the following definitions:

**Definition 2.17.** Given random elements  $X_1, \dots, X_n$  on the same underlying probability space,  $\mathcal{L}(X_1, \dots, X_n) := \mathcal{L}((X_1, \dots, X_n))$  is called the *joint distribution* of  $X_1, \dots, X_n$ . Conversely, given a random tuple  $(X_1, \dots, X_n)$ , each  $\mathcal{L}(X_i)$  is called a *marginal distribution*.

Suppose we have two distributions of random elements  $\mathcal{L}(X_1)$  and  $\mathcal{L}(X_2)$ . *Coupling* is the technique of constructing a random ordered pair  $(X_1, X_2)$  which realizes the given distributions as marginal distributions. Usually this is done by specifying the joint distribution  $\mathcal{L}(X_1, X_2)$ .

The idea is that coupling creates a particular kind of dependence between  $X_1$  and  $X_2$  that allows us to compare the two distributions. Often, we are able to make conclusions about the distributions  $\mathcal{L}(X_i)$  which are independent of their specific realizations as random elements in the coupling.

## 2.3 Markov Chains

**Comment 2.2.** I'll need to define Markov Chains, stationary distributions, irreducibility and time-reversibility.

Perhaps I should talk more generally about stochastic processes, because applying exchangeable pairs to Stein's method has connections with Ornstein-Uhlenbeck processes and also Stein's method can be applied to Poisson processes.

## 2.4 The Weak Topology on Probability Measures

**Comment 2.3.** The main purpose of this section is to motivate the metrics usually used in Stein's method: they are all legitimate topological metrics and are consistent with the topology of convergence in distribution. In particular, if we can show  $d_{\mathcal{H}}(X_n, X) \rightarrow 0$  we have shown that  $X_n \xrightarrow{d} X$ , as Toby does [Joh11].

**Definition 2.18.** Let  $(X_n)_{n \in \mathbb{N}}$  be a sequence of random variables. We say  $X_n$  *converges in distribution* to a random variable  $X$  if  $\mathbb{E}f(X_n) \rightarrow \mathbb{E}f(X)$  for all bounded continuous functions  $f$ . Alternatively, we say  $\mathcal{L}(X_n)$  converges *weakly* to  $\mathcal{L}(X)$ , or simply  $\mathcal{L}(X_n) \rightarrow \mathcal{L}(X)$ . The topology on  $\mathcal{P}(\mathbb{R})$  associated with this convergence is called the *weak topology* (we will see that it is indeed a topology). Convergence in distribution of random vectors is defined component-wise.

**Definition 2.19.** The *distribution function*  $F_X$  of a random variable  $X$  is defined by  $F_X(x) = \mathbb{P}(X \leq x)$ .

**Theorem 2.20.** *The following are equivalent.*

1.  $\mathcal{L}(X_n) \rightarrow \mathcal{L}(X)$
2.  $F_{X_n}(x) \rightarrow F_X(x)$  for all  $x$  where  $F_X$  is continuous
3. (*Lévy's continuity theorem*)  $\mathbb{E}e^{itX_n} \rightarrow \mathbb{E}e^{itX}$  for all  $t \in \mathbb{R}$ .

The equivalence of Items 1 and 2 is a well-known result called the Portmanteau Theorem.

When  $X$  and each  $X_n$  are integer random variables, then Item 2 reduces to the condition that  $\mathbb{P}(X_n = k) \rightarrow \mathbb{P}(X = k)$  for all  $k$ . This characterization is usually used to prove the Poisson limit theorem. Lévy's continuity theorem is classically used to prove the central limit theorem, but we will not discuss it in this thesis.



For combinatorial applications, convergence in distribution can also be proved by the “method of moments”: if  $X$  is the only random variable with the moments  $(\mathbb{E}X^k)_{k \in \mathbb{N}}$ , then  $\mathcal{L}(X_n) \rightarrow \mathcal{L}(X)$  if  $\mathbb{E}X_n^k \rightarrow \mathbb{E}X^k$ . Convergence in distribution can also sometimes be inferred from stronger forms of convergence when  $X$  and all the  $X_n$  are coupled to the same underlying space.

A disadvantage of all these approaches is that it is difficult to quantify the rate of convergence.

In functional analysis terms, note that expectation operators are bounded linear functionals on the space of real bounded continuous functions. Then,  $\mathcal{L}(X_n) \rightarrow \mathcal{L}(X)$  just means that  $\mathbb{E}_{X_n} \rightarrow \mathbb{E}_X$  in the weak-star topology. Although  $C_b(\mathbb{R})^*$  is not metrizable, the subspace corresponding to  $\mathcal{P}(\mathbb{R})$  is in fact metrizable, with a metric called the Lévy metric. For Stein’s method we will be interested in some slightly stronger metrics.

**Definition 2.21.** Let  $\mathcal{H}$  be a collection of real measurable “test” functions. Define  $d_{\mathcal{H}} : \mathcal{P}(\mathbb{R})^2 \rightarrow \mathbb{R}^+$  by  $d_{\mathcal{H}}(\mathbb{P}_1, \mathbb{P}_2) = \sup_{h \in \mathcal{H}} |\mathbb{E}_{\mathbb{P}_1} h - \mathbb{E}_{\mathbb{P}_2} h|$ . For random variables  $X_1, X_2$ , we write  $d_{\mathcal{H}}(X_1, X_2)$  instead of  $d_{\mathcal{H}}(\mathcal{L}(X_1), \mathcal{L}(X_2))$ .

**Issue 2.4.** What if there is no supremum or  $\mathbb{E}_{\mathbb{P}_1} h = \infty$ ?

Each  $d_{\mathcal{H}}$  is non-negative, symmetric and satisfies the triangle inequality.

**Definition 2.22.** A set of real functions  $\mathcal{H}$  is a *determining class* if  $\mathbb{E}_{\mathbb{P}_1} h = \mathbb{E}_{\mathbb{P}_2} h$  for all  $h \in \mathcal{H}$  implies that  $\mathbb{P}_1 = \mathbb{P}_2$ .

To check that  $d_{\mathcal{H}}$  is a metric, we only need to check that  $d_{\mathcal{H}}(\mathbb{P}_1, \mathbb{P}_2) = 0$  implies that  $\mathbb{P}_1 = \mathbb{P}_2$ . That is, we need to check that  $\mathcal{H}$  is a determining class.

**Definition 2.23.** We define some special cases of  $d_{\mathcal{H}}$ .

- If  $\mathcal{H}_K = \{1_{(-\infty, x]} : x \in \mathbb{R}\}$  then  $d_K := d_{\mathcal{H}_K}$  is called the *Kolmogorov metric*.

- If  $\mathcal{H}_W$  is the set of real functions  $h$  that satisfy  $|h(x_1) - h(x_2)| \leq |x_1 - x_2|$  for all  $x_1, x_2 \in \mathbb{R}$  (that is, the set of functions with Lipschitz constant 1), then  $d_W := d_{\mathcal{H}_W}$  is called the *Wasserstein metric*.
- If  $\mathcal{H}_{TV}$  is the set of functions  $1_B$  for Borel  $B$ ,  $d_{TV} := d_{\mathcal{H}_{TV}}$  is called the *total variation metric*.

**Proposition 2.24.** *The Kolmogorov, Wasserstein and total variation “metrics” are actually metrics.*

*Proof.* We check that  $\mathcal{H}_K$ ,  $\mathcal{H}_W$  and  $\mathcal{H}_{TV}$  are determining classes. Let  $\mathcal{H} \in \{\mathcal{H}_K, \mathcal{H}_W, \mathcal{H}_{TV}\}$ , and suppose that  $\mathbb{E}_{\mathbb{P}_1} h = \mathbb{E}_{\mathbb{P}_2} h$  for all  $h \in \mathcal{H}$ . It suffices to prove that  $\mathbb{P}_1((-\infty, x]) = \mathbb{P}_2((-\infty, x])$  for all  $x \in \mathbb{R}$ , since the sets  $(-\infty, x]$  generate the Borel  $\sigma$ -algebra. For  $\mathcal{H} \in \{\mathcal{H}_K, \mathcal{H}_{TV}\}$  this is immediate, because  $\mathbb{P}_i((-\infty, x]) = \mathbb{E}_{\mathbb{P}_i} 1_{(-\infty, x]}$ . So, consider,  $\mathcal{H} = \mathcal{H}_W$ .

For  $\varepsilon > 0$  and  $x \in \mathbb{R}$ , let  $h_{x,\varepsilon}$  be the continuous function which takes the value 1 on the set  $(-\infty, x]$ , takes the value 0 on the set  $[x + \varepsilon, \infty)$ , and is linearly interpolated in the range  $[x, x + \varepsilon]$ . Since  $\varepsilon h_{x,\varepsilon} \in \mathcal{H}_W$ , we have  $\mathbb{E}_{\mathbb{P}_1} h_{x,\varepsilon} = \mathbb{E}_{\mathbb{P}_2} h_{x,\varepsilon}$  for each  $n \in \mathbb{N}$ . For each  $x \in \mathbb{R}$ ,  $h_{x,1/n} \rightarrow 1_{(-\infty, x]}$  pointwise and each  $h_{x,1/n} \leq 1$  so by the dominated convergence theorem,  $\mathbb{E}_{\mathbb{P}_i} h_{1/n} \rightarrow \mathbb{E}_{\mathbb{P}_i} 1_{(-\infty, x]}$  for each  $i \in \{1, 2\}$ . We have again proved that  $\mathbb{P}_1((-\infty, x]) = \mathbb{P}_2((-\infty, x])$  for all  $x \in \mathbb{R}$ .  $\square$

**Proposition 2.25.** *The topologies induced by the Kolmogorov, Wasserstein and total variation metrics are each stronger than the weak topology.*

*Proof.* If  $d_K(X_n, X) \rightarrow 0$  or  $d_{TV}(X_n, X) \rightarrow 0$  then  $F_{X_n} \rightarrow F_X$  uniformly, so certainly Item 2 of Theorem 2.20 holds.

Now, suppose  $d_K(X_n, X) \rightarrow 0$ . Let  $d_n = \sqrt{d_K(X_n, X)}$  and recall the definition of  $h_{x,\varepsilon}$  from the proof of Proposition 2.24. Since  $d_n h_{x,d_n} \in \mathcal{H}_K$  for each  $n \in \mathbb{N}$ , we have  $\mathbb{E}_{X_n} h_{x,d_n} - \mathbb{E}_X h_{x,d_n} \leq d_K(X_n, X)/d_n = d_n \rightarrow 0$  uniformly for  $x \in \mathbb{R}$ . Now, note that  $F_X(x - \varepsilon) \leq \mathbb{E}_X h_{x-\varepsilon,\varepsilon} \leq$

$F_X(x) \leq \mathbb{E}_X h_{x,\varepsilon} \leq F_X(x + \varepsilon)$  for any random variable  $X$ . If  $F_X$  is continuous at  $x$  then

$$F_{X_n}(x) - F_X(x) \leq (\mathbb{E}_{X_n} h_{x,d_n} - \mathbb{E}_X h_{x,d_n}) + (F_X(x + d_n) - F_X(x)) \rightarrow 0$$

$$F_{X_n}(x) - F_X(x) \geq (\mathbb{E}_{X_n} h_{x-d_n,d_n} - \mathbb{E}_X h_{x-d_n,d_n}) + (F_X(x - d_n) - F_X(x)) \rightarrow 0$$

so Item 2 of Theorem 2.20 holds. □

Proposition 2.25 tells us that we can sensibly use our metrics to quantify the distance between random variables, in a way that is consistent with distributional (weak) convergence. All three metrics are relevant in their own right, but sometimes one may be easier to work with. It is sometimes possible to transfer results between metrics, though this usually results in worse constants than working directly in the desired metric.

**Issue 2.5.** It may be worthwhile to actually characterize the Wasserstein, Kolmogorov and Total Variation topologies. In particular, Wikipedia says that Wasserstein convergence is just weak convergence plus convergence of the first moment.

**Definition 2.26.** If  $F_X(x) = \int_{-\infty}^x f_X(x) dx$  then  $f_X$  is called the *Lebesgue density* of  $X$ , and  $X$  is called a *continuous* random variable.

If  $X$  is a continuous random variable, then by the Radon-Nikodym chain rule  $\mathbb{E}_X h = \int_{\mathbb{R}} h(x) f_X(x) dx$ .

**Proposition 2.27.** *Let  $X_1, X_2$  be random variables.*

1.  $d_K(X_1, X_2) \leq d_{TV}(X_1, X_2)$
2. *If  $X_2$  has Lebesgue density bounded by  $C$ , then  $d_K(X_1, X_2) \leq \sqrt{2C d_W(X_1, X_2)}$ .*

*Proof.* (Adapted from [Ros11, Proposition 1.2]). Item 1 is immediate from the definition.

Then, as in the proof of Proposition 2.25,

$$\begin{aligned}
F_{X_n}(x) - F_X(x) &\leq (\mathbb{E}_{X_n} h_{x,\varepsilon} - \mathbb{E}_X h_{x,\varepsilon}) + (\mathbb{E}_X h_{x,\varepsilon} - F_X(x)) \\
&\leq d_W(X_1, X_2)/\varepsilon + \int_x^{x+\varepsilon} h_{x,\varepsilon} f_X(x) \, dx \\
&\leq d_W(X_1, X_2)/\varepsilon + C\varepsilon/2
\end{aligned}$$

and similarly

$$F_{X_n}(x) - F_X(x) \geq -d_W(X_1, X_2)/\varepsilon - C\varepsilon/2,$$

So, we can take  $\varepsilon = \sqrt{2d_W(X_1, X_2)/C}$  to prove Item 2.  $\square$

**Example 2.28.** If  $\mathcal{L}_{X_2} = \mathcal{N}(0, 1)$  then  $d_K \leq (2/\pi)^{1/4} \sqrt{d_W(X_1, X_2)}$ .

In a combinatorial setting, many of our results are about integer random variables. The total variation metric is usually exclusively used in this case.

**Proposition 2.29.** *If  $X_1, X_2$  are integer-valued random variables, then*

$$d_{TV}(X_1, X_2) = \frac{1}{2} \sum_{k \in \mathbb{Z}} |\mathbb{P}(X_1 = k) - \mathbb{P}(X_2 = k)|.$$

*Proof.* For any Borel set  $A$ , let  $d_A = \mathbb{P}(X_1 \in A) - \mathbb{P}(X_2 \in A)$ , so that  $d_{TV}(X_1, X_2) = \sup |d_A|$ . Define

$$A_{<} = \{k \in \mathbb{Z} : \mathbb{P}(X_1 = k) < \mathbb{P}(X_2 = k)\},$$

$$A_{>} = \{k \in \mathbb{Z} : \mathbb{P}(X_1 = k) > \mathbb{P}(X_2 = k)\}.$$

For any Borel  $A$ , we have

$$\begin{aligned}
d_A &= \sum_{k \in \mathbb{Z} \cap A} (\mathbb{P}(X_1 = k) - \mathbb{P}(X_2 = k)) \\
&\leq \sum_{k \in A_{>}} (\mathbb{P}(X_1 = k) - \mathbb{P}(X_2 = k)) \\
&= d_{A_{>}}
\end{aligned}$$

and similarly  $\mathbb{P}(X_1 \in A) - \mathbb{P}(X_2 \in A) \geq d_{A_<}$ . Since  $d_{A_>} = -d_{A_<}$ , we have

$$d_{\text{TV}}(X_1, X_2) = (d_{A_>} - d_{A_<})/2 = \frac{1}{2} \sum_{k \in \mathbb{Z}} |\mathbb{P}(X_1 = k) - \mathbb{P}(X_2 = k)|.$$

□

### 3 Random Combinatorial Structures

**Definition 3.1.** Given a finite space of combinatorial objects  $\Omega$ , a probability space  $(\Omega, 2^\Omega, \mathbb{P})$  is often called a *model* of  $\Omega$ .

**Definition 3.2.** In a probability space  $(\Omega, 2^\Omega, \mathbb{P})$  where  $\Omega$  is finite, if  $\mathbb{P}(\omega) = 1/|\Omega|$  for each  $\omega \in \Omega$ , then we say the space is *uniform*.

Uniform models are the simplest examples of random structures. For example, the uniform space  $\mathcal{S}_n$  of permutations on  $n$  elements has  $\mathbb{P}(\sigma) = 1/n!$  for each  $\sigma \in \mathcal{S}_n$ . The uniform random graph model  $\mathcal{G}_{n,M}$  has  $\mathbb{P}(G) = \binom{\binom{n}{2}}{M}^{-1}$  for each graph  $G$  on the vertex set  $[n]$  which has  $M$  edges. The uniform random regular graph model  $\mathcal{G}_{n,d}$  is uniform on the set of all  $d$ -regular graphs on the vertex set  $[n]$ , though an explicit formula for the number of such graphs is not known.

As an important example of a (generally) non-uniform model, the (Erdős-Rényi) binomial random graph model  $\mathcal{G}_{n,p}$  has

$$\mathbb{P}(G) = p^{|E(G)|} (1-p)^{\binom{n}{2} - |E(G)|}$$

for each graph  $G$  on the vertex set  $[n]$ . When  $p = 1/2$ , we obtain the uniform model on all graphs on the vertex set  $[n]$ .

One way to conceptualize the binomial model is to consider a sequence of independent coin tosses, where the coin is biased to land heads with probability  $p$ . Each coin toss corresponds to a particular potential edge, and determines whether that edge is present in the final random

graph. When we define more complicated random models, we will often use this kind of informal description rather than giving an explicit formula for each  $\mathbb{P}(\omega)$ .

As another example, the uniform model  $\mathcal{G}_{n,M}$  can be alternatively defined recursively:  $\mathcal{G}_{n,0}$  is always the trivial graph with no edges, and for each  $M > 0$ , to obtain  $\mathcal{G}_{n,M}$  we choose  $G \in \mathcal{G}_{n,M-1}$  and add one of the  $\binom{n}{2} - (M-1)$  possible edges at random.

**Comment 3.1.** This section is unfinished, I'll probably want random matrices and maybe the pairing model on random regular graphs

## 4 Stein's Method in the Abstract

**Comment 4.1.** There are a few quite different presentations of Stein's method. One thing I'm trying to do here is to unify Stein's functional analysis approach for exchangeable pairs [Ste86] with Ross' general presentation [Ros11].

The reason I want to look at Stein's original, more abstract presentation is that I think it does a better job motivating why things work. Before I read that, the steps taken to apply Stein's method seemed like blindly doing things and it turns out they work.

Suppose we have a potentially complicated random variable  $X$ , and we believe the distribution of  $X$  is close to a “standard” distribution  $\mathcal{L}_0$ . Then, Stein's method allows us to compare the operators  $\mathbb{E}_X$  and  $\mathbb{E}_0 := \mathbb{E}_{\mathcal{L}_0}$ . This is sometimes directly useful for approximating statistics of  $X$  (for example,  $\mathbb{P}(X \in A) = \mathbb{E}_X 1_A$ ). However, particularly for combinatorial applications, Stein's method is most often used to bound the distance  $d_{\mathcal{H}}(\mathcal{L}_X, \mathcal{L}_0)$ , where the metric  $d_{\mathcal{H}}$  from Definition 2.21 is defined in terms of  $\mathbb{E}_X$  and  $\mathbb{E}_0$ .

Stein's method is motivated by the idea of a characterizing operator.

**Definition 4.1.** Let  $\mathcal{F}_0$  be a vector space and  $\mathcal{X}_0$  be a vector space of measurable functions. We say a linear operator  $T_0 : \mathcal{F}_0 \rightarrow \mathcal{X}_0$  is a *characterizing operator* for the distribution  $\mathcal{L}_0$  if  $\text{im } T_0 = \mathcal{X}_0 \cap \ker \mathbb{E}_0$ . For convenience, where there is no ambiguity we will often implicitly restrict  $\mathbb{E}_0$  to  $\mathcal{X}_0$ , so we can write  $\text{im } T_0 = \ker \mathbb{E}_0$ .

The following proposition shows why  $T_0$  is called a characterizing operator.

**Proposition 4.2.** *If  $T_0 : \mathcal{F}_0 \rightarrow \mathcal{X}_0$  is a characterizing operator and  $\mathcal{X}_0$  is a determining class then  $\text{im } T_0 \subseteq \ker \mathbb{E}_X$  implies  $\mathcal{L}_X = \mathcal{L}_0$ .*

*Proof.* If  $h \in \mathcal{X}_0$ , then  $h - \mathbb{E}_0 h \in \ker \mathbb{E}_0 = \text{im } T_0$  so  $\mathbb{E}_X[h - \mathbb{E}_0 h] = 0$ . That is,  $\mathbb{E}_X h = \mathbb{E}_0 h$  for all  $h \in \mathcal{X}_0$ , which means  $\mathcal{L}_X = \mathcal{L}_0$  by the definition of a determining class.  $\square$

**Issue 4.2.** Ross [Ros11] and others use this weaker condition as the definition of a characterizing operator. I'll have to look at examples of operators that satisfy the weaker but not the stronger condition to see if the stronger definition is warranted (my guess is yes, if Stein decided to originally define it the way I did).

**Proposition 4.3.**  *$T_0 : \mathcal{F}_0 \rightarrow \mathcal{X}_0$  is characterizing if and only if there is a linear operator  $U_0 : \mathcal{X}_0 \rightarrow \mathcal{F}_0$  such that the following two equations hold.*

$$\mathbb{E}_0 T_0 = 0_{\mathcal{F}_0}, \tag{4.1}$$

$$T_0 U_0 + \mathbb{E}_0 = \text{id}_{\mathcal{X}_0}. \tag{4.2}$$

*Proof.* Suppose  $T_0$  is a characterizing operator. Equation (4.1) is immediate. Let  $\{h_i\}_{i \in \mathcal{I}}$  be a (Hamel) basis of  $\mathcal{X}_0$ . For each  $i \in \mathcal{I}$  we have  $h_i - \mathbb{E}_0 h_i \in \ker \mathbb{E}_0$  so there is some  $f_i$  (not necessarily unique) that solves  $T_0 f_i = h_i - \mathbb{E}_0 h_i$ . The operator  $U_0$  can then be defined by  $\sum_{i \in \mathcal{I}} a_i h_i \mapsto \sum_{i \in \mathcal{I}} a_i f_i$ , satisfying (4.2).

**Issue 4.3.** there's probably a cleaner functional analysis way to prove that. Also, is  $U_0$  bounded?

Conversely, suppose (4.1) holds and  $U_0$  exists satisfying (4.2). For  $h \in \ker \mathbb{E}_0$  we have  $T_0(U_0h) = h$  and  $h \in \operatorname{im} T_0$ , so  $\ker \mathbb{E}_0 \subseteq \operatorname{im} T_0$ . Equation (4.1) immediately says that  $\operatorname{im} T_0 \subseteq \ker \mathbb{E}_0$ , so  $T_0$  is a characterizing operator.  $\square$

We'll use Proposition 4.3 to give two important examples of characterizing operators.

**Theorem 4.4.** Define  $T_{\mathcal{N}}$  by  $T_{\mathcal{N}}f(x) = f'(x) - xf(x)$ . Let  $\mathcal{X}_{\mathcal{N}} = L^1(\mathbb{R}, \mathcal{N}(0, 1))$  be the set of functions  $h : \mathbb{R} \rightarrow \mathbb{R}$  that satisfy  $\mathbb{E}_{\mathcal{N}}|h| < \infty$  and let  $\mathcal{F}_{\mathcal{N}} = T_{\mathcal{N}}^{-1}\mathcal{X}_{\mathcal{N}}$  be the set of functions  $f : \mathbb{R} \rightarrow \mathbb{R}$  such that  $\mathbb{E}_{\mathcal{N}}|T_{\mathcal{N}}f| < \infty$ . Then  $T_{\mathcal{N}} : \mathcal{F}_{\mathcal{N}} \rightarrow \mathcal{X}_{\mathcal{N}}$  is a characterizing operator for  $\mathcal{N}(0, 1)$ .

*Proof.* For any  $f \in \mathcal{F}_{\mathcal{N}}$ , integration by parts gives

$$\mathbb{E}_{\mathcal{N}}T_{\mathcal{N}}f = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-t^2/2} f'(t) - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} te^{-t^2/2} f(t) \, dt = 0$$

so  $\mathbb{E}_{\mathcal{N}}T_{\mathcal{N}}f = 0$  and (4.1) holds. Then, define  $U_{\mathcal{N}}$  by

$$U_{\mathcal{N}}h(x) = e^{x^2/2} \int_{-\infty}^x (h(t) - \mathbb{E}_{\mathcal{N}}h) e^{-t^2/2} \, dt.$$

By the product rule and the fundamental theorem of calculus, for all  $h \in \mathcal{X}_{\mathcal{N}}$  we have

$$T_{\mathcal{N}}U_{\mathcal{N}}h(x) = h(x) - \mathbb{E}_{\mathcal{N}}h,$$

so (4.2) holds and Proposition 4.3 completes the proof.  $\square$

**Theorem 4.5.** Define  $T_{\operatorname{Po}(\lambda)}$  by  $T_{\operatorname{Po}(\lambda)}f(k) = \lambda f(k+1) - kf(k)$ . Let  $\mathcal{X}_{\operatorname{Po}(\lambda)}$  be the set of integer-valued functions  $h : \mathbb{N} \rightarrow \mathbb{Z}$  that satisfy  $\mathbb{E}_{\operatorname{Po}(\lambda)}|h| < \infty$  and let  $\mathcal{F}_{\operatorname{Po}(\lambda)} = T_{\operatorname{Po}(\lambda)}^{-1}\mathcal{X}_{\operatorname{Po}(\lambda)}$  be the set of functions  $f : \mathbb{N} \rightarrow \mathbb{Z}$  such that  $\mathbb{E}_{\operatorname{Po}(\lambda)}|T_{\operatorname{Po}(\lambda)}f| < \infty$ . Then  $T_{\operatorname{Po}(\lambda)} : \mathcal{F}_{\operatorname{Po}(\lambda)} \rightarrow \mathcal{X}_{\operatorname{Po}(\lambda)}$  is a characterizing operator for  $\operatorname{Po}(\lambda)$ .



*Proof.* For any  $f \in \mathcal{F}_{\text{Po}(\lambda)}$ , we have

$$\mathbb{E}_{\text{Po}(\lambda)} T_{\text{Po}(\lambda)} f = e^{-\lambda} \sum_{i=0}^{\infty} \frac{\lambda^{i+1}}{i!} f(i+1) - e^{-\lambda} \sum_{i=1}^{\infty} \frac{\lambda^i}{(i-1)!} f(i) = 0$$

so  $\mathbb{E}_{\text{Po}(\lambda)} T_{\text{Po}(\lambda)} = 0$  and (4.1) holds. Then, define  $U_{\text{Po}(\lambda)}$  by

$$U_{\text{Po}(\lambda)} h(k) = \frac{(k-1)!}{\lambda^k} \sum_{i=0}^{k-1} \frac{\lambda^i}{i!} (h(i) - \mathbb{E}_{\text{Po}(\lambda)} h)$$

for  $k \geq 1$ . Substituting and simplifying gives

$$T_{\text{Po}(\lambda)} U_{\text{Po}(\lambda)} h(k) = h(k) - \mathbb{E}_{\text{Po}(\lambda)} h,$$

so (4.2) holds and Proposition 4.3 completes the proof.  $\square$

Note that  $\mathcal{H}_{\text{TV}} \subseteq \mathcal{X}_{\mathcal{N}}$ , where  $\mathcal{H}_{\text{TV}}$  is as defined in Definition 2.23. Since  $\mathcal{H}_{\text{TV}}$  is a determining class,  $T_{\mathcal{N}}$  is a characterizing operator in the sense of Proposition 4.2. We can say the same about  $T_{\text{Po}(\lambda)}$  if we restrict our attention to integer-valued random variables.

**Issue 4.4.** Stein chose  $\mathcal{X}_0 = \{h : \mathbb{E}[\text{id}_{\mathbb{R}}^k |h|] < \infty \text{ for all } k\}$ , for both the Poisson and normal case. I'm not sure why, I'll revisit this after looking at exchangeable pairs.

The utility of the introduction of a characterizing operator is that for each  $h \in \mathcal{X}_0$ , Equation (4.2) allows us to make the transformation

$$\mathbb{E}_X h = \mathbb{E}_0 h + \mathbb{E}_X T_0 U_0 h. \tag{4.3}$$

The original purpose of Stein's method was to estimate some particular  $\mathbb{E}_X h$ . If  $\mathcal{L}_0$  was chosen to be a “simple”, well-understood distribution then the term  $\mathbb{E}_0 h$  should be easy to compute or estimate, and if the distribution of  $X$  was “close” to  $\mathcal{L}_0$ , then it should be possible to show that the remainder  $\mathbb{E}_X T_0 U_0 h$  is small.

For our purposes, the main use of (4.3) is to bound  $d_{\mathcal{H}}(X, \mathcal{L}_0)$  for some  $\mathcal{H} \subseteq \mathcal{X}_0$ . For any  $\mathcal{Y} \supseteq U_0\mathcal{H}$ , we have

$$d_{\mathcal{H}}(X, \mathcal{L}_0) = \sup_{h \in \mathcal{H}} |\mathbb{E}_X T_0 U_0 h| \leq \sup_{f \in \mathcal{Y}} |\mathbb{E}_X T_0 f|.$$

We have reduced the problem of bounding  $d_{\mathcal{H}}(X, \mathcal{L}_0)$  to that of bounding  $|\mathbb{E}_X T_0 f|$  (uniformly over  $f \in \mathcal{Y}$ ). Especially in the cases where  $\mathcal{L}_0$  is normal or Poisson and  $\mathcal{H}$  is one of the standard choices in Definition 2.23, there are a number of known convenient choices of  $\mathcal{Y}$ , and a number of methods that are known to be effective to bound  $|\mathbb{E}_X T_0 f|$ .

**Example 4.6.** If  $\mathcal{H} = \mathcal{H}_{\text{TV}}$  and  $\mathcal{L}_0 = \text{Po}(\lambda)$ , using the characterizing operator in Theorem 4.5, then we generally choose

$$\mathcal{Y} = \left\{ f \in \mathcal{F}_0 : \|f\|_{\infty} \leq \min\{1, \lambda^{-1/2}\}, \|\Delta f\|_{\infty} \leq \min\{1, \lambda^{-1}\} \right\},$$

where  $\Delta f(k) = f(k+1) - f(k)$ .

Proving that this choice of  $\mathcal{Y}$  satisfies  $\mathcal{Y} \supseteq U_{\text{Po}(\lambda)}\mathcal{H}$  is nontrivial. But, we can prove that the constraints are of the “correct” order of magnitude.

**Issue 4.5.** The proof is in [BHJ92, Remark 10.2.4]. There’s also a simpler proof in [BHJ92, Lemma 1.1.1] that  $\|f\|_{\infty} \leq 2 \min\{1, \lambda^{-1/2}\}$  suffices. I’ll revisit this later.

**Proposition 4.7.** *With  $U_{\text{Po}(\lambda)}$  as in Theorem 4.5, we have*

$$\begin{aligned} \sup_A \|U_{\text{Po}(\lambda)} 1_A\|_{\infty} &\asymp \lambda^{-1/2} \text{ as } \lambda \rightarrow \infty, & \sup_A \|U_{\text{Po}(\lambda)} 1_A\|_{\infty} &\asymp 1 \text{ as } \lambda \rightarrow 0, \\ \sup_A \|\Delta U_{\text{Po}(\lambda)} 1_A\|_{\infty} &\asymp \lambda^{-1} \text{ as } \lambda \rightarrow \infty, & \sup_A \|\Delta U_{\text{Po}(\lambda)} 1_A\|_{\infty} &\asymp 1 \text{ as } \lambda \rightarrow 0, \end{aligned}$$

*Proof.* (Adapted from [BHJ92, Lemma 1.1.1]) For any Borel  $A$  and any  $k \in \mathbb{N}$ ,

$$\begin{aligned} U_{\text{Po}(\lambda)} 1_A(k) &= \frac{(k-1)!}{\lambda^k} (\text{Po}_\lambda(A \cap [k-1]) - \text{Po}_\lambda(A) \text{Po}_\lambda([k-1])) \\ &= \frac{(k-1)!}{\lambda^k} \left( \text{Po}_\lambda(A \cap [k-1]) (1 - \text{Po}_\lambda([k-1])) \right. \\ &\quad \left. - (\text{Po}_\lambda(A) - \text{Po}_\lambda(A \cap [k-1])) \text{Po}_\lambda([k-1]) \right) \\ &= \frac{(k-1)!}{\lambda^k} \left( \text{Po}_\lambda(A \cap [k-1]) \text{Po}_\lambda(\mathbb{R} \setminus [k-1]) - \text{Po}_\lambda(A \setminus [k-1]) \text{Po}_\lambda([k-1]) \right). \end{aligned}$$

Note that  $\text{Po}_\lambda(A \cap [k-1])$  is bounded above by  $\text{Po}_\lambda([k-1])$  and  $\text{Po}_\lambda(A \setminus [k-1])$  is bounded above by  $\text{Po}_\lambda(\mathbb{R} \setminus [k-1])$ , so

$$|U_{\text{Po}(\lambda)} 1_A(k)| \leq \frac{(k-1)!}{\lambda^k} \text{Po}_\lambda([k-1]) \text{Po}_\lambda(\mathbb{R} \setminus [k-1]).$$

Note that we have equality when  $A = [k-1]$ . If  $k \asymp \lambda$  as  $\lambda \rightarrow \infty$  then Stirling's approximation (unfinished...) □

**Issue 4.6.** I need a way to show  $\text{Po}_\lambda([k-1]) \asymp 1$  when  $k$  is close to  $\lambda$ . I imagine the typical approach would be to show that  $(\text{Po}(\lambda) - \lambda)/\lambda \rightarrow \mathcal{N}(0, 1)$  with Lévy's continuity theorem, but that would require me to introduce more stuff in the probability revision section. Maybe I can use Stein's method itself for this! In any case I'll think about this more when I've decided on my approach to Issue 4.5.

**Issue 4.7.** I should have a toy example here that is amenable to several methods.

## 4.1 The method of Exchangeable Pairs

This is Stein's original approach, and is effective in wide generality for discrete random variables. In what follows, we assume  $X$  is discrete. Let  $\Omega_X$  be the support of  $X$ .

**Example 4.8** (adapted from [Ros11, Example 4.21]). We will use an example problem to illustrate the principles in this section. We will say a *fixed point* of a permutation  $\sigma \in S_n$  is an index  $k \in [n]$  that satisfies  $\sigma(k) = k$ . Let  $X : S_n \rightarrow \{0\} \cup [n]$  give the number of fixed points in each permutation from  $S_n$ . We interpret  $X$  as a random variable on the underlying space  $S_n$ .

Now, if  $n$  is large then fixed points are largely independent of each other, and each of  $n$  indices has a probability of  $1/n$  to be a fixed point. So (recalling Remark 2.10), we might expect  $\mathcal{L}_X$  to be “close” to  $\text{Po}(1)$ . We will attempt to bound  $d_{\text{TV}}(X, \text{Po}(1))$  to quantify this intuition.

The general idea is that we can use an object  $\mathbf{X}$  called an exchangeable pair to construct a characterizing operator  $T_{\mathbf{X}}$  for  $X$ . We then use an operator  $\alpha$  to connect the domains of  $T_{\mathbf{X}}$  and  $T_0$  in such a way that  $T_{\mathbf{X}}\alpha$  approximates  $T_0$ . We then have  $\mathbb{E}_X T_0 = \mathbb{E}_X (T_0 - T_{\mathbf{X}}\alpha)$ , so we can use the fact that  $T_0 - T_{\mathbf{X}}\alpha$  is small to bound  $\mathbb{E}_X T_0 f$ .

**Definition 4.9.** A 2-dimensional random pair  $\mathbf{X} = (X_1, X_2)$  is an *exchangeable pair* if  $\mathcal{L}(X_1, X_2) = \mathcal{L}(X_2, X_1)$ . We will denote the support of  $\mathbf{X}$  by  $\Omega_{\mathbf{X}}^{(2)}$ .

That is, a pair  $\mathbf{X}$  is exchangeable if exchanging the components of the pair does not change their joint distribution. In particular, the marginal distribution of  $X_1$  and  $X_2$  must be the same.

**Comment 4.8.** All presentations of Stein’s method I’ve seen use the notation  $(X, X')$  but I think that has the potential to be confusing because the  $X$  in that pair is defined on  $\Omega^2$  whereas the original random variable  $X$  is defined on  $\Omega$ .

**Proposition 4.10.** *There is a natural equivalence between time-homogeneous reversible Markov chains with steady-state distribution  $\mathcal{L}_X$ , and exchangeable pairs with margins  $\mathcal{L}_X$ .*

*Proof.* Given an exchangeable pair  $\mathbf{X}$  with margins  $\mathcal{L}_X$ , we can define a time-homogeneous Markov chain  $M$  with transition probabilities  $p(x_1, x_2) = \mathbb{P}(X_2 = x_2 | X_1 = x_1)$ . With  $\pi(x) =$

$\mathbb{P}(X = x)$ , we then have

$$\pi(x_1)p(x_1, x_2) = \mathbb{P}(\mathbf{X} = (x_1, x_2)) = \pi(x_2)p(x_2, x_1)$$

for any  $x_1, x_2 \in \Omega_X$ . So,  $M$  is reversible with steady-state distribution  $\mathcal{L}_X$ .

Conversely, suppose we have a time-homogeneous reversible Markov chain with steady-state distribution  $\mathcal{L}_X$ . Let the transition probability between  $x$  and  $x'$  be  $p(x, x')$  and let the probability of state  $x$  in the steady-state distribution be  $\pi(x)$ . We can then define an exchangeable pair  $\mathbf{X}$  by

$$\mathbb{P}(\mathbf{X} = (x_1, x_2)) = \pi(x_1)p(x_1, x_2) = \pi(x_2)p(x_2, x_1) = \mathbb{P}(\mathbf{X} = (x_2, x_1))$$

and the proposition is proved.  $\square$

**Definition 4.11.** We say an exchangeable pair  $\mathbf{X}$  is *connected* if its Markov chain is irreducible.

*Remark 4.12.* We are particularly interested in exchangeable pairs  $\mathbf{X}$  with marginal distributions  $\mathcal{L}_{X_1} = \mathcal{L}_{X_2} = \mathcal{L}_X$ . If  $X$  is defined on an underlying combinatorial probability space  $(\Omega, 2^\Omega, \mathbb{P})$ , it is often convenient to first construct an exchangeable pair  $\mathbf{W} = (W_1, W_2)$  with margins  $\mathbb{P}$ , so that the vector  $\mathbf{X}_{\mathbf{W}} = (X(W_1), X(W_2))$  is an exchangeable pair with margins  $\mathcal{L}_X$ . If  $(W_1, W_2)$  is connected, then  $\mathbf{X}_{\mathbf{W}}$  is connected also.

**Example 4.13.** We continue Example 4.8. We will define a specific exchangeable pair  $\mathbf{W} = (W_1, W_2)$  with margins  $\mathcal{S}_n$  by

$$\mathbb{P}((W_1, W_2) = (\sigma_1, \sigma_2)) = \begin{cases} (n! \binom{n}{2})^{-1} & \text{if } \sigma_1 = \sigma_2(ij) \text{ for some transposition } (ij) \\ 0 & \text{otherwise.} \end{cases}$$

The relation of differing by a transposition is symmetric, so  $\mathbf{W}$  is indeed an exchangeable pair. The Markov chain associated with  $\mathbf{W}$  has a simple interpretation. Given a random permutation  $\sigma$ , to make a transition in the Markov chain we just randomly choose one of the  $\binom{n}{2}$  possible transpositions and compose it with  $\sigma$ . Because the transpositions generate

$S_n$ , the pair  $\mathbf{W}$  is connected, so we can use the construction from Remark 4.12 to produce a connected exchangeable pair  $\mathbf{X}$  with margins  $\mathbf{W}$ .

The Markov chain underlying a connected exchangeable pair can be naturally viewed as a connected one-dimensional simplicial complex. The zeroth reduced homology group  $\ker \partial_0 / \text{im } \partial_1$  of a connected simplicial complex has dimension zero, and this motivates the construction of a characterizing operator in a natural way. (The following theorem is self-contained and requires no knowledge of homology theory).

**Theorem 4.14.** *Suppose  $\mathbf{X}$  is a connected exchangeable pair with margins  $\mathcal{L}_X$ . Let  $\mathcal{F}_X \subseteq L^1(\Omega_X^2, \mathcal{L}_X)$  be the set of functions  $f : \Omega_X^2 \rightarrow \mathbb{R}$  that are antisymmetric in the sense that  $f(x_1, x_2) = -f(x_2, x_1)$  and satisfy  $\mathbb{E}|f(\mathbf{X})| < \infty$ . Let  $\mathcal{X}_X = L^1(\Omega_X, \mathcal{L}_X)$  be the set of functions  $h : \Omega_X \rightarrow \mathbb{R}$  that satisfy  $\mathbb{E}_X|h| < \infty$ .*

*Define  $T_{\mathbf{X}} : \mathcal{F}_X \rightarrow \mathcal{X}_X$  by  $T_{\mathbf{X}}f(x) = \sum_{x' \in \Omega_X} f(x, x_2)p(x, x_2) = \mathbb{E}[f(\mathbf{X})|X_1 = x]$ , so that  $T_{\mathbf{X}}X = \mathbb{E}^{X_1}f(\mathbf{X})$ . Then  $T_{\mathbf{X}}$  is a characterizing operator for  $X$ .*

*Proof.* To see that  $\text{im } T_{\mathbf{X}} \subseteq \ker \mathbb{E}_X$ , fix  $f \in \mathcal{F}$  and note that by the tower law of expectation (Proposition 2.16),

$$\mathbb{E}_X T_{\mathbf{X}}f = \mathbb{E} \mathbb{E}^{X_1} f(\mathbf{X}) = \mathbb{E} f(\mathbf{X}).$$

By exchangeability and antisymmetry,  $\mathbb{E}f(\mathbf{X}) = \mathbb{E}f(X_2, X_1) = -\mathbb{E}f(\mathbf{X})$ , so  $\mathbb{E}f(X_1, X_2) = \mathbb{E}_X T_{\mathbf{X}}f = 0$ . This did not require the connectedness condition. We can similarly prove that  $T_{\mathbf{X}}$  is well-defined as an operator from  $\mathcal{F}_X$  to  $\mathcal{X}_X$ : note that  $\mathbb{E}_X|T_{\mathbf{X}}f| = \mathbb{E}|f(\mathbf{X})|$  so  $T_{\mathbf{X}}\mathcal{F}_X \subseteq \mathcal{X}_X$ .

We will next prove  $\ker \mathbb{E}_X \subseteq \text{im } T_{\mathbf{X}}$ , but first we make some definitions. For each  $x \in \Omega_X$ , let  $h_x$  be the function that takes the value  $\pi(x)^{-1}$  on  $x$  and is zero elsewhere, so that  $h = \sum_{x \in \Omega_X} h(x)\pi(x)h_x$  for each  $h \in \mathcal{X}_X$ . For each  $(x_1, x_2) \in \Omega_{\mathbf{X}}^{(2)}$ , define  $f_{x_1, x_2} \in \mathcal{F}_X$  as the function that takes the value  $(\pi(x_1)p(x_1, x_2))^{-1}$  on  $(x_1, x_2)$ , takes the value  $-(\pi(x_2)p(x_2, x_1))^{-1}$  on  $(x_2, x_1)$ , and takes the value zero elsewhere. Note that this function is antisymmetric by the reversibility of the Markov chain of  $\mathbf{X}$ . We have  $T_{\mathbf{X}}f_{x_1, x_2} = h_{x_2} - h_{x_1}$ .

Let  $h \in \ker \mathbb{E}_X$ , and fix arbitrary  $x^* \in \Omega_X$ . By the connectedness assumption, for each  $x \in \Omega_X$  there is a sequence

$$x = x^{(0)}, x^{(1)}, \dots, x^{(k-1)}, x^{(k)} = x^*$$

with  $(x^{(i-1)}, x^{(i)}) \in \Omega_{\mathbf{X}}^{(2)}$  for  $i = 1, \dots, k$ . Note that

$$h_{x^*} - h_x = T_{\mathbf{X}} \sum_{i=1}^k f_{x^{(i-1)}, x^{(i)}} =: T_{\mathbf{X}} f_x^*,$$

and it follows that

$$h = \sum_{x \in \Omega_X} h(x) \pi(x) (h_{x^*} - T_{\mathbf{X}} f_x^*) = (\mathbb{E}_X h) h_{x^*} - \sum_{x \in \Omega_X} T_{\mathbf{X}} h(x) \pi(x) f_x^*.$$

By assumption  $(\mathbb{E}_X h) = 0$ . If  $\Omega_X$  is finite, as it will be in our applications, then it would immediately follow that  $h \in \text{im } T_{\mathbf{X}}$ . Otherwise we will need some functional analysis. Note that  $T_{\mathbf{X}}$  is actually an isometry between a subspace  $\mathcal{F}_X$  of  $L^1(\Omega_X^2, \mathcal{L}_{\mathbf{X}})$  and  $L^1(\Omega_X, \mathcal{L}_X)$ . First we prove that  $\mathcal{F}_X$  is closed and therefore a Banach space. For any  $f \in L^1(\Omega_X^2, \mathcal{L}_{\mathbf{X}})$ , let  $\bar{f}$  be defined by  $(x_1, x_2) \mapsto -f(x_2, x_1)$ , so that  $f \in \mathcal{F}_X$  implies that  $f = \bar{f}$ . Suppose  $f_n \rightarrow f$ , with  $f_n \in \mathcal{F}_X$  for all  $n \in \mathbb{N}$ . By exchangeability we have

$$\|f_n - f\| = \mathbb{E}|f_n(\mathbf{X}) - f(\mathbf{X})| = \mathbb{E}|f_n(X_2, X_1) - \bar{f}(X_2, X_1)| = \|f_n - \bar{f}\|$$

so  $f = \bar{f}$  and  $f \in \mathcal{F}_X$ . Finally, it is a simple fact that an isometry between Banach spaces has a closed range. For, if  $(T_{\mathbf{X}} f_n)_{n \in \mathbb{N}}$  is a Cauchy sequence in  $\text{im } T_{\mathbf{X}}$ , then  $(f_n)_{n \in \mathbb{N}}$  is a Cauchy sequence in  $\mathcal{F}_X$  and converges to some  $f \in \mathcal{F}_X$ . It follows that  $T_{\mathbf{X}} f_n \rightarrow T_{\mathbf{X}} f \in \text{im } T_{\mathbf{X}}$ .

We have proved that  $h \in \text{im } T_{\mathbf{X}}$ , completing the proof that  $\ker \mathbb{E}_X \subseteq \text{im } T_{\mathbf{X}}$ .  $\square$

The final step is to choose an operator  $\alpha : \mathcal{F}_0 \rightarrow \mathcal{F}_X$  in such a way that  $T_{\mathbf{X}}$  can be easily compared with  $T_0 \alpha$ .

**Example 4.15.** For the Poisson case in Theorem 4.5, we need to compare  $\lambda f(X+1)$  with

$Xf(X)$ . It is often fruitful to define  $\alpha$  by

$$\alpha f(x_1, x_2) = cf(x_2) 1\{x_2 = x_1 + 1\} - cf(x_1) 1\{x_1 = x_2 + 1\}$$

for some  $c \in \mathbb{R}$ . We will then have

$$T_0 f - T_{\mathbf{X}} \alpha f = f(X + 1)(1 - c\mathbb{P}(X_2 = X_1 + 1|X_1)) - f(X_1)(X_1 - c\mathbb{P}(X_1 = X_2 + 1|X_1)).$$

Using the triangle inequality and the choice of  $\mathcal{Y}$  in Example 4.6, for all  $f \in \mathcal{Y}$ :

$$\mathbb{E}_X T_0 f \leq \min\left(1, \lambda^{-1/2}\right)(\mathbb{E}|1 - c\mathbb{P}(X_2 = X_1 + 1|X_1)| + \mathbb{E}|X_1 - c\mathbb{P}(X_1 = X_2 + 1|X_1)|).$$

This approximation is effective when

$$\mathbb{P}(X_2 = X_1 + 1|X_1) \approx \lambda/c, \tag{4.4}$$

$$\mathbb{P}(X_1 = X_2 + 1|X_1) \approx X_1/c.$$

The interpretation of these approximate equalities is that the Markov chain associated with  $\mathbf{X}$  is approximately an immigration-death process. This is likely to happen when  $X(\omega)$  is in some sense a statistic of the amount of local structure over the object  $\omega$ , and  $\mathbf{X}$  is defined by a Markov chain on  $\Omega$  (as in Remark 4.12) that (uniformly) randomly disturbs local structure. The conclusion to Example 4.8 should make this clear:

**Example 4.16.** We continue Example 4.8, recalling the exchangeable pairs  $\mathbf{W}$  and  $\mathbf{X}$  from Example 4.13. The interpretation of

$$\mathbb{P}(X_1 = X_2 + 1|X_1) = \mathbb{P}(X_2 = X_1 - 1|X_1)$$

is the probability of a transposition destroying exactly one out of an existing  $X_1$  fixed points. In order to destroy exactly one fixed point, we have to choose a fixed point to destroy, and swap it with a non-fixed-point. There are  $X_1(n - X_1)$  out of  $\binom{n}{2}$  transpositions that do this, so

$$\mathbb{P}(X_1 = X_2 + 1|X_1) = \frac{X_1(n - X_1)}{\binom{n}{2}}.$$



Next, we will find a formula for  $\mathbb{P}(X(W_2) = X(W_1) + 1|W_1)$ , noting that

$$\mathbb{P}(X_2 = X_1 + 1|X_1) = \mathbb{E}[\mathbb{P}(X(W_2) = X(W_1) + 1|W_1)|X_1].$$

In order to create exactly one fixed point, we have to choose an index  $k$  that is not fixed in  $W_1$  (there are  $n - X_1$  such) and compose  $W_1$  with  $(k \sigma(k))$ . This creates exactly one fixed point unless  $\sigma^{-1}(k) = k$ , in which case it creates two. We have counted this second case twice for every transposition in the cycle decomposition of  $W_1$ . Let  $Y$  be the number of transpositions in the cycle decomposition of  $W_1$ . We have

$$\mathbb{P}(X_2 = X_1 + 1|X_1) = \frac{n - X_1 - 2\mathbb{E}[Y|X_1]}{\binom{n}{2}}.$$

In order to satisfy (4.4) as closely as possible, we choose  $c = \binom{n}{2}/n$ . Recalling that  $\mathbb{E}X_1 = 1$ , we then have

$$\begin{aligned} \mathbb{E}|1 - c\mathbb{P}(X_2 = X_1 + 1|X_1)| &= \mathbb{E}\left[1 - \frac{n - X_1 - 2\mathbb{E}[Y|X_1]}{n}\right] \\ &= 1/n + 2\mathbb{E}Y/n, \\ \mathbb{E}|X_1 - c\mathbb{P}(X_2 = X_1 - 1|X_1)| &= \mathbb{E}\left[X_1 - \frac{X_1(n - X_1)}{n}\right] \\ &= \mathbb{E}[X_1^2]/n. \end{aligned}$$

Now, the probability that a transposition  $(i j)$  is in the cycle decomposition of  $W_1$  is  $(n(n - 1))^{-1}$  because  $i$  must map to  $j$  out of the  $n$  possible options in  $[n]$ , then  $j$  must map to  $i$  out of the  $n - 1$  possible options in  $[n] \setminus \{j\}$ . There are  $\binom{n}{2} = n(n - 1)/2$  possible transpositions so by Remark 2.10 it follows that  $\mathbb{E}Y = 1/2$ .

Now,  $\mathbb{E}[X_1(X_1 - 1)/2]$  is the expected number of unordered pairs of distinct fixed points in a permutation  $\sigma \in \mathcal{S}_n$ . For any unordered pair of distinct indices  $\{i, j\}$ , the probability that both are fixed is  $(n(n - 1))^{-1}$  because  $i$  must map to  $i$  out of the  $n$  possible options in  $[n]$ , then  $j$  must map to  $j$  out of the  $n - 1$  possible options in  $[n] \setminus \{i\}$ . The total number of unordered pairs is  $\binom{n}{2} = n(n - 1)/2$ , so again applying Remark 2.10, we have  $\mathbb{E}[X_1(X_1 - 1)/2] = 1/2$  and  $\mathbb{E}[X_1^2] = 2$ .

We conclude that  $d_{\text{TV}}(X, \text{Po}(1)) \leq 4/n$ .

**Comment 4.9.** The “generator method” [BC05] says that the Poisson characterizing operator can be obtained with the generator of an immigration-death process and the Normal characterizing operator can be obtained with the generator of an O-U process. Investigate the link here?

## 4.2 Size-Bias Coupling

## Part II

# Applications

**Comment 4.10.** I’d like to go into a number of small examples (perhaps interspersed in the discussion of Stein’s method in Part I), but I’d like to also go through a number of “big” examples. I’d like these examples to showcase

- different types of results: most applications give quantitative estimates. [Joh11] gives a non-quantitative distributional convergence result that was not previously proved using other methods. There are also results that have no connection with distribution metrics, such as the concentration inequalities in [Ros11]. In particular, the Latin rectangle example in [Ste86] is interesting in that the final result is not probabilistic.
- different types of distributions: definitely at least the Poisson and normal case, perhaps also an example of a more exotic distribution like the one in [FG12] or perturbations of Poisson/normal distributions as in [BČX07].
- different ways to apply Stein’s method: definitely exchangeable pairs and probably size-biasing. Maybe also Zero-bias coupling.

## References

- [BC05] Andrew D Barbour and Louis Hsiao Yun Chen, *An introduction to Stein's method*, vol. 4, World Scientific, 2005.
- [BČX07] Andrew D Barbour, Vydas Čekanavičius, and Aihua Xia, *On Stein's method and perturbations*, arXiv preprint math/0702008 (2007).
- [BHJ92] Andrew D Barbour, Lars Holst, and Svante Janson, *Poisson approximation*, Clarendon press Oxford, 1992.
- [FG12] Jason Fulman and Larry Goldstein, *Stein's method and the rank distribution of random matrices over finite fields*, arXiv preprint arXiv:1211.0504 (2012).
- [Joh11] Tobias Johnson, *Exchangeable pairs, switchings, and random regular graphs*, arXiv preprint arXiv:1112.0704 (2011).
- [Ros11] Nathan Ross, *Fundamentals of Stein's method*, Probab. Surv **8** (2011), 210–293.
- [Ste86] Charles Stein, *Approximate computation of expectations*, Lecture Notes – Monograph Series **7** (1986).